

# Virtualization and Cloud Computing at Fermilab

## NLIT Summit 2011

Keith Chadwick

Grid & Cloud Computing Department  
Fermilab

Work supported by the U.S. Department of Energy under contract No. DE-AC02-07CH11359

# Acknowledgements

- Many of these slides have been copied from a recent Fermilab Computing Division “briefing” on Virtualization and Cloud Computing.
- The contributors to that briefing include:
  - Grid & Cloud Computing Department,
  - Fermilab Experiments Facilities Department,
  - Virtual Services Group,
  - Services Operations Support Department,
  - Stakeholders (Grid, CMS, REX, CET, OSG, etc.).

# Outline

- FermiGrid Initial Strategy and History
- FermiGrid-HA Project and Virtualization
- FermiGrid-HA2 Project
- GridWorks
- FermiCloud
- Other Virtualization Projects
  - FEF VM Clusters & General Physics Compute Facility (GPCF)
  - Virtual Services
- Some Performance Measurements
- Summary and Final Thoughts
  - Workloads
  - Open Source vs. Commercial
  - Xen vs. KVM
  - Security

# FermiGrid – Initial Strategy

- Strategy:
  - In order to better serve the entire program of Fermilab, the Computing Division has undertaken the strategy of placing all of its production resources in a Grid "meta-facility" infrastructure called FermiGrid.
- This strategy is designed to allow Fermilab:
  - to insure that the large experiments who currently have dedicated resources to have first priority usage of those resources that are purchased on their behalf.
  - to allow opportunistic use of these dedicated resources, as well as other shared Farm and Analysis resources, by various Virtual Organizations (VO's) that participate in the Fermilab experimental program and by certain VOs that use the Open Science Grid (OSG).
  - to optimise use of resources at Fermilab.
  - to make a coherent way of putting Fermilab on the Open Science Grid.
  - to save some effort and resources by implementing certain shared services and approaches.
  - to work together more coherently to move all of our applications and services to run on the Grid.
  - to better handle a transition from Run II to LHC in a time of shrinking budgets and possibly shrinking resources for Run II worldwide.
  - to fully support Open Science Grid and the LHC Computing Grid and gain positive benefit from this emerging infrastructure in the US and Europe.



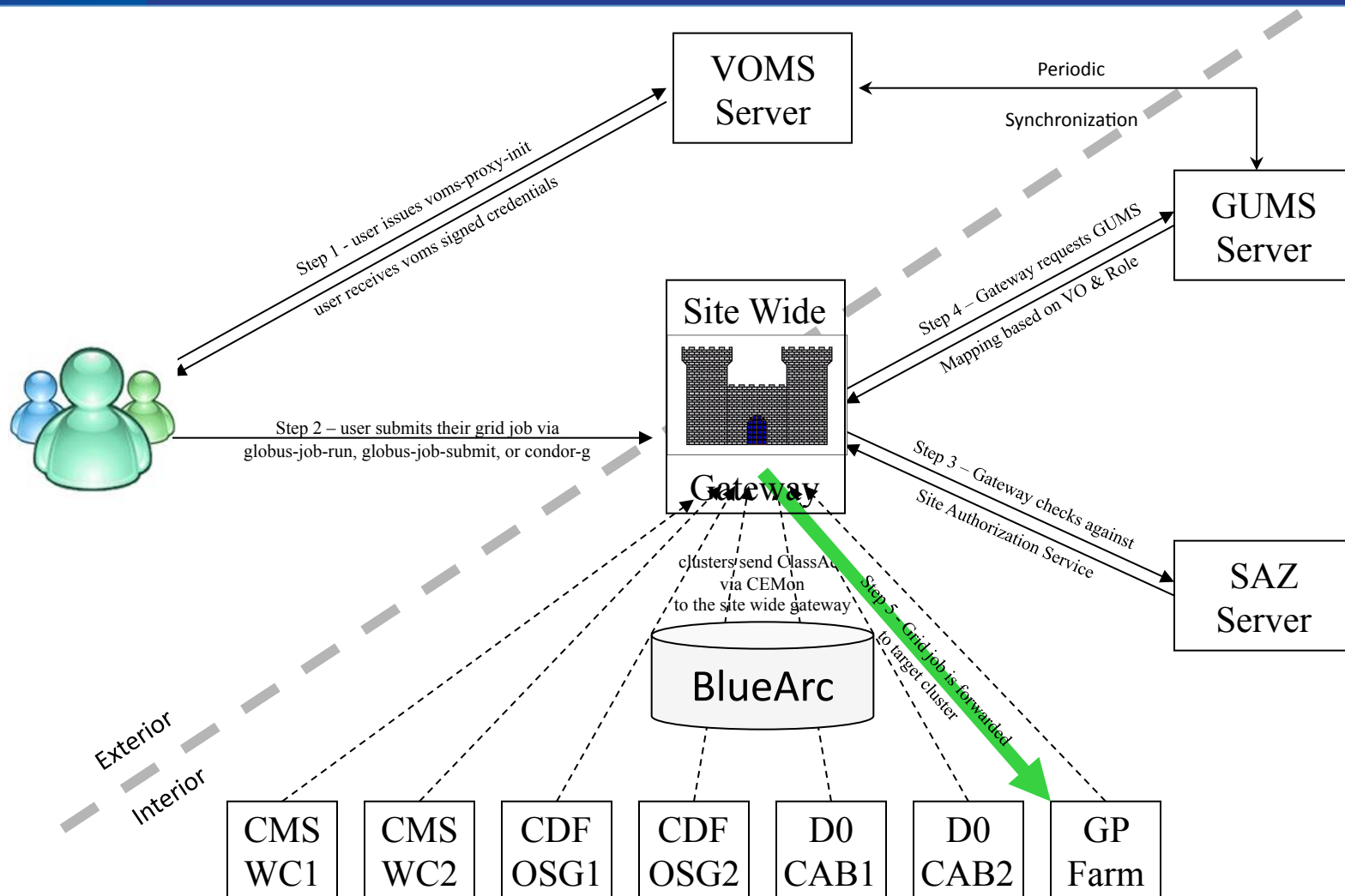
# What is FermiGrid?

- FermiGrid is:
  - The Fermilab campus Grid and Grid portal.
    - The site globus gateway.
    - Accepts jobs from external (to Fermilab) sources and forwards the jobs onto internal clusters.
  - A set of common services to support the campus Grid and interface to Open Science Grid (OSG) / LHC Computing Grid (LCG):
    - VOMS, VOMRS, GUMS, SAZ, MyProxy, Squid, Gratia Accounting, etc.
  - A forum for promoting stakeholder interoperability and resource sharing within Fermilab:
    - CMS, CDF, D0;
    - ktev, miniboone, minos, mipp, etc.
  - The Open Science Grid portal to Fermilab Compute and Storage Services.
- FermiGrid Web Site & Additional Documentation:
  - <http://fermigrid.fnal.gov/>

# FermiGrid – The Early Years

- The FermiGrid concept started in mid CY2004 and the FermiGrid strategy was formally announced in late CY2004.
  - The initial hardware was ordered and delivered in early CY2005.
- The Initial core services (Globus Gateway, VOMS and GUMS) based on OSG 0.2.1 were commissioned on April 1, 2005.
  - We missed the ides of March, so we chose April Fools Day...
- Our first site gateway which used Condor-G matchmaking and MyProxy was commissioned in the fall of CY2005.
  - Job forwarding based on work by GridX1 in Canada (<http://www.gridx1.ca>).
  - Users were required to store a copy of their delegated grid proxy in our MyProxy repository prior to using the job forwarding gateway.
- OSG 0.4 was deployed across FermiGrid in late January-February 2006.
  - Followed quickly by OSG 0.4.1 in March 2006.
- The Site AuthoriZation (SAZ) service was commissioned on October 2, 2006.
  - Provides site wide whitelist and blacklist capability.
  - Can make decision based on any of DN, VO, Role, and CA.
  - Currently operate in a default accept mode (providing that the presented proxy was generated using voms-proxy-init).
- The glexec pilot job glide-in service was commissioned on November 1, 2006.
  - Provides authorization and accounting trail for Condor glide-in jobs.
- An upgraded version of the site job forwarding gateway (jobmanager-cemon) was commissioned in November 2006
  - Eliminated the need to utilize MyProxy via "accept limited" option on the gatekeeper.
  - Based on CEMon and OSG RESS, Condor Matchmaking.
  - Periodic hold and release functions were added in March 2007.
- OSG 0.6.0 was deployed across FermiGrid during March & April 2007.

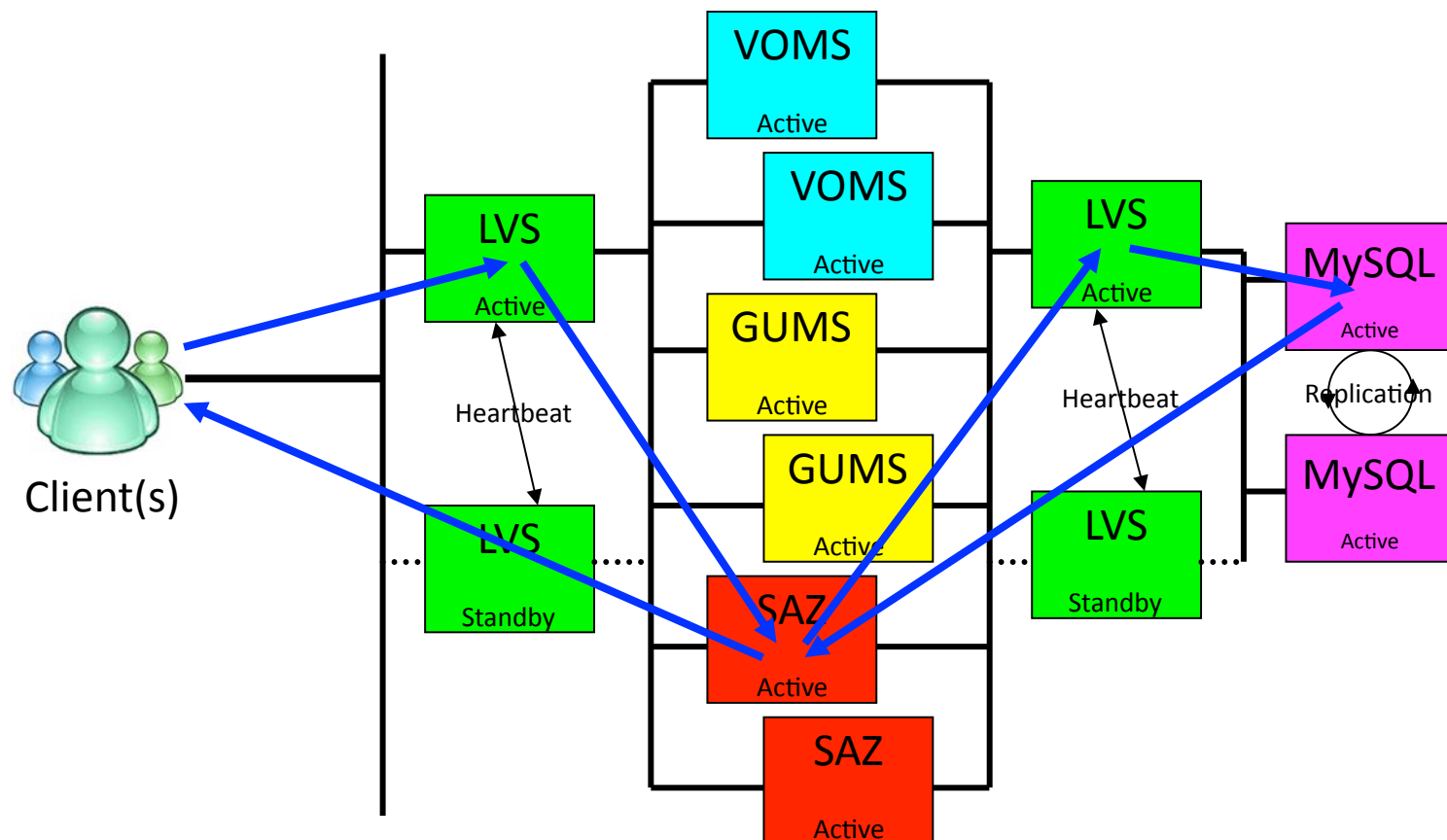
# Initial FermiGrid Architecture



# FermiGrid-HA Project

- The FermiGrid service deployment was working reasonably well and meeting the negotiated service level agreements, but we recognized that an outage of a critical service (GUMS, SAZ) would result in the entire Grid at Fermilab “going dark”.
- The FermiGrid-HA (High Availability) project was established to review the FermiGrid service design and make the necessary changes to support highly available services.
- We designed a redundant service architecture that could be deployed as a “drop in” replacement for the current services.

# FermiGrid HA Services - 1



# So what about Virtualization?

- Unfortunately, the initial design of the “drop in” replacement redundant service architecture required a significant increase in the hardware to host the services”.
- We did not have sufficient budget to buy all of the necessary hardware.
- Simultaneously, we had been performing virtualization explorations using open source Xen during 2006 to mid 2007.
- The decision was made to incorporate virtualization to the FermiGrid-HA project plan in order to reduce the required hardware footprint (ultimately to just two systems).

# FermiGrid-HA Services - 2

## Xen Domain 0

LVS      Xen VM 0  
Active      fg5x0

VOMS      Xen VM 1  
Active      fg5x1

GUMS      Xen VM 2  
Active      fg5x2

SAZ      Xen VM 3  
Active      fg5x3

MySQL      Xen VM 4  
Active      fg5x4

Active

fermigrd5

## Xen Domain 0

LVS      Xen VM 0  
Standby      fg6x0

VOMS      Xen VM 1  
Active      fg6x1

GUMS      Xen VM 2  
Active      fg6x2

SAZ      Xen VM 3  
Active      fg6x3

MySQL      Xen VM 4  
Active      fg6x4

Active

fermigrd6

# FermiGrid-HA – Key Technologies

- FermiGrid-HA utilizes the following key technologies:
  - Scientific Linux 5 + Xen Hypervisor.
  - Linux Virtual Server (LVS) in a Direct Routing (DR) configuration.
  - MySQL Circular Replication.
  - DRBD Volume Replication (for MyProxy).



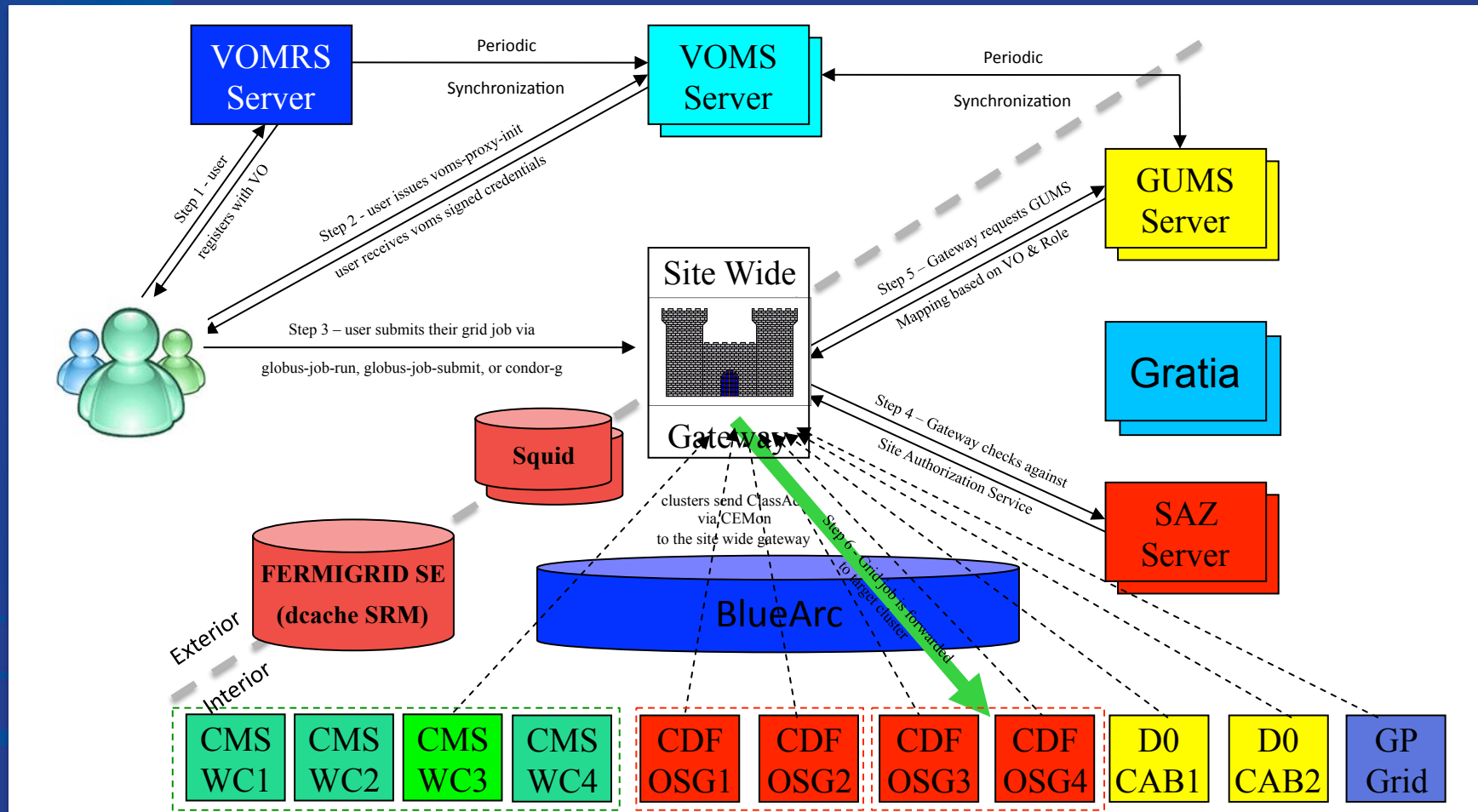
# FermiGrid employs several strategies to deploy HA services

- Trivial monitoring or information services (such as Ganglia and Zabbix) are deployed on two independent virtual machines.
- Services that natively support HA operation (Condor Information Gatherer, FermiGrid internal ReSS deployment) are deployed in the standard service HA configuration on two independent virtual machines.
- Services that maintain intermediate routing information (Linux Virtual Server) are deployed in an active/passive configuration on two independent virtual machines. A periodic heartbeat process is used to perform any necessary service failover.
- Services that do not maintain intermediate context (i.e. are pure request/response services such as GUMS, SAZ, Squid) are deployed using a Linux Virtual Server (LVS) front end to active/active servers on two independent virtual machines.
- Services that support active-active database functions (circularly replicating MySQL servers) are deployed on two independent virtual machines.
- Services that require “real-time” filesystem synchronization (MyProxy) are deployed using DRBD in an active-passive configuration on two independent virtual machines.

# FermiGrid-HA – Virtualization and Highly Available Grid Services

- The FermiGrid core services deployment was virtualized using paravirtualized Xen as part of FermiGrid-HA deployment in the fall of 2007.
- Virtualized remaining services using Xen during 2008.
  - Statically deployed cloud computing.
- The goal for FermiGrid-HA is > 99.999% service availability.
  - Not including Building or Network failures.
  - FermiGrid actively measures the service availability of the services in the FermiGrid service catalog.
  - <http://fermigrd.fnal.gov/fermigrd-metrics.html>
- High availability has been demonstrated:
  - For the first seven months, we achieved a service availability of 99.9969%.
  - 100% availability for critical services (VOMS, GUMS, SAZ, Squid) for more than a year.
  - Redundancy has also allowed rolling upgrades & patches without service interruption.
- High performance MySQL has been demonstrated:
  - OSG and Fermilab Gratia accounting databases.
- We have just completed the FermiGrid-HA2 project.
  - Redundant services hosted in two buildings (FCC and GCC).
  - No longer have to exclude building failures in availability measurements.

# Current FermiGrid Architecture



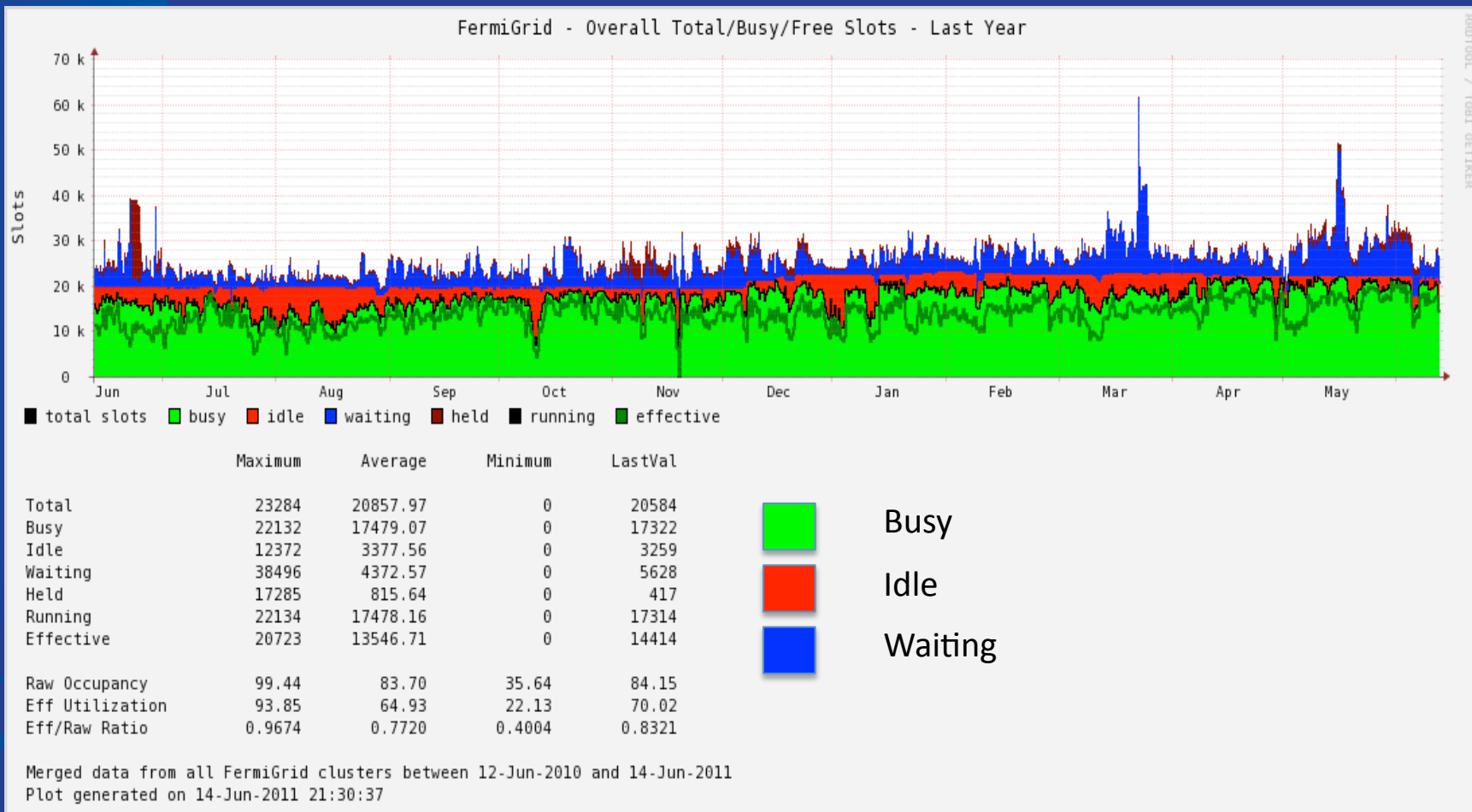
# Measured FermiGrid Service Availability for the Past Year\*

Service	Availability	Downtime
VOMS	100%	0m
GUMS	100%	0m
SAZ (gatekeeper)	100%	0m
Squid	100%	0m
MyProxy	99.943%	297.6m
ReSS	99.957%	223.7m
Gratia	99.892%	565.5m

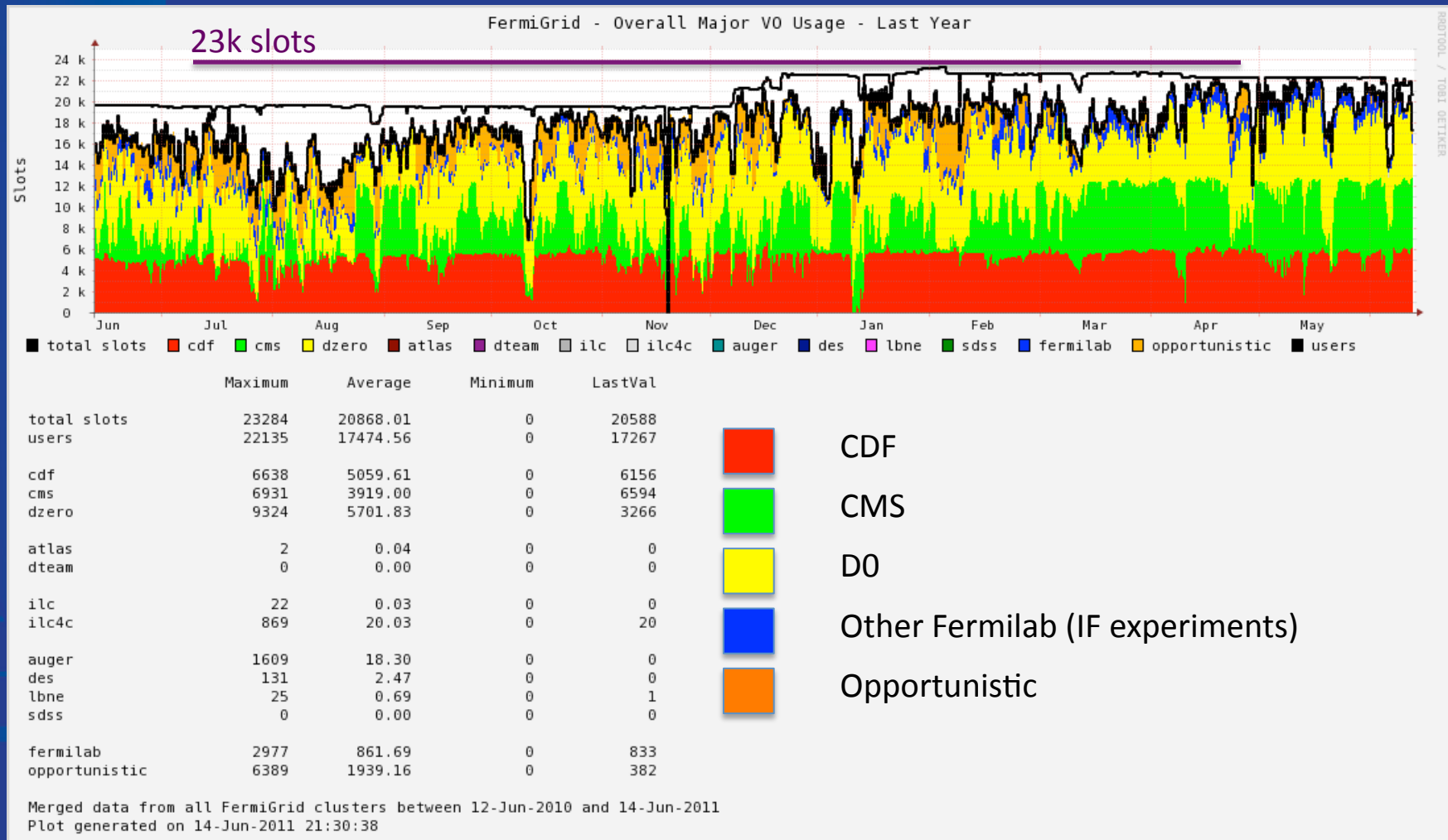
\* = Excluding building or network failures and scheduled downtimes

# FermiGrid

## Overall Occupancy & Utilization



# FermiGrid – Past Year Slot Usage



# Batch Systems, Occupancy & Utilization

Cluster	Cluster Batch System	Current Cluster Size (Slots)	Average Cluster Occupancy (%)
CDF	Condor	5,600	89.3
CMS	Condor	7,485	88.8
D0	PBS	6,916	73.4
GP	Condor	3,284	68.7
Overall	----	23,285	82.3

# SAZ – Central Banning Service

- Site AuthoriZation service, developed at Fermilab,
- Allows us to ban any DN, VO, FQAN, or CA,
- Don't have to wait for CA to revoke the certificate or VO to remove from membership,
- Available as Pacman package,
- Details at <http://saz.fnal.gov>
- With glideins, every worker node can be a client simultaneously (have seen 5000+ active clients),
- Significant work has been done to make it more resilient and scalable,
- Code also contains support for new XACML-based authorization protocol.

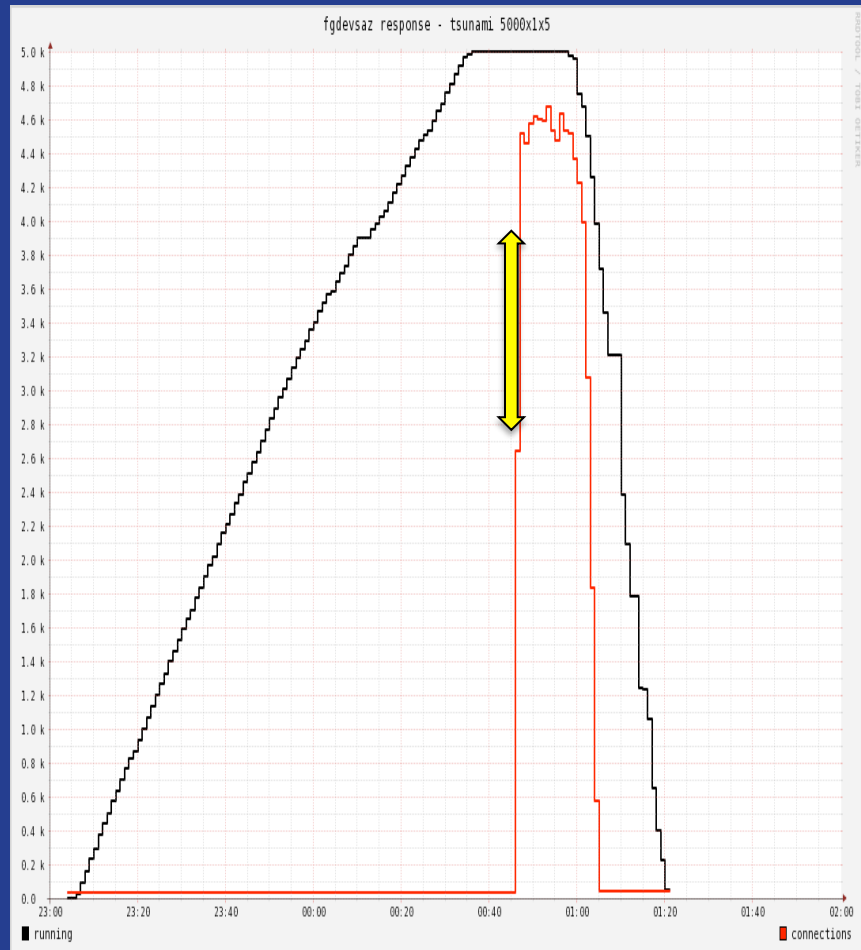


# Why a new SAZ Server?

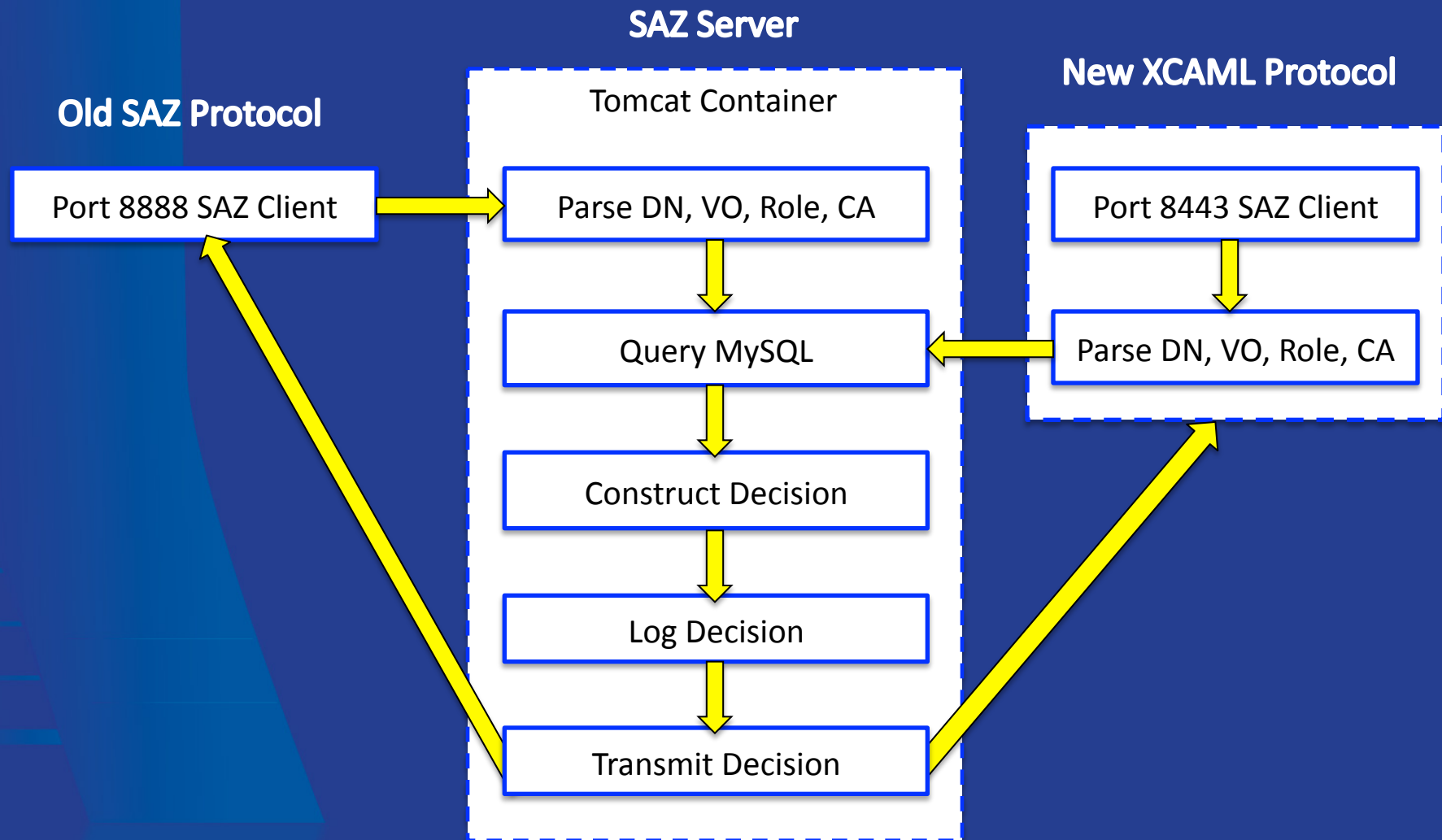
- Previous SAZ server (V2.0.1b) had shown itself extremely vulnerable to user generated authorization “tsunamis”:
  - Very short duration jobs
  - User issues condor\_rm on a large (>1000) glidein.
  - We had replicated multiple SAZ servers (on virtual machines), but users could still overwhelm the service.
- This was fixed in the new SAZ Server (V2.7.2)
  - using tomcat and pools of execution and hibernate threads.
  - Various other bugs were found and fixed in the current SAZ server and sazclient.
  - Added support for the XACML protocol (used by Globus).
  - NOT transitioning to using XACML (yet).
- New SAZ Servers (V2.7.2) have been in operation since May 2010.
  - Incrementally deployed/upgraded in the existing virtual machines.
  - Deployed without any downtime.
  - Have repeatedly demonstrated in production that they are robust against authorization “tsunamis”.

# SAZ Server V2.7.2 Performance

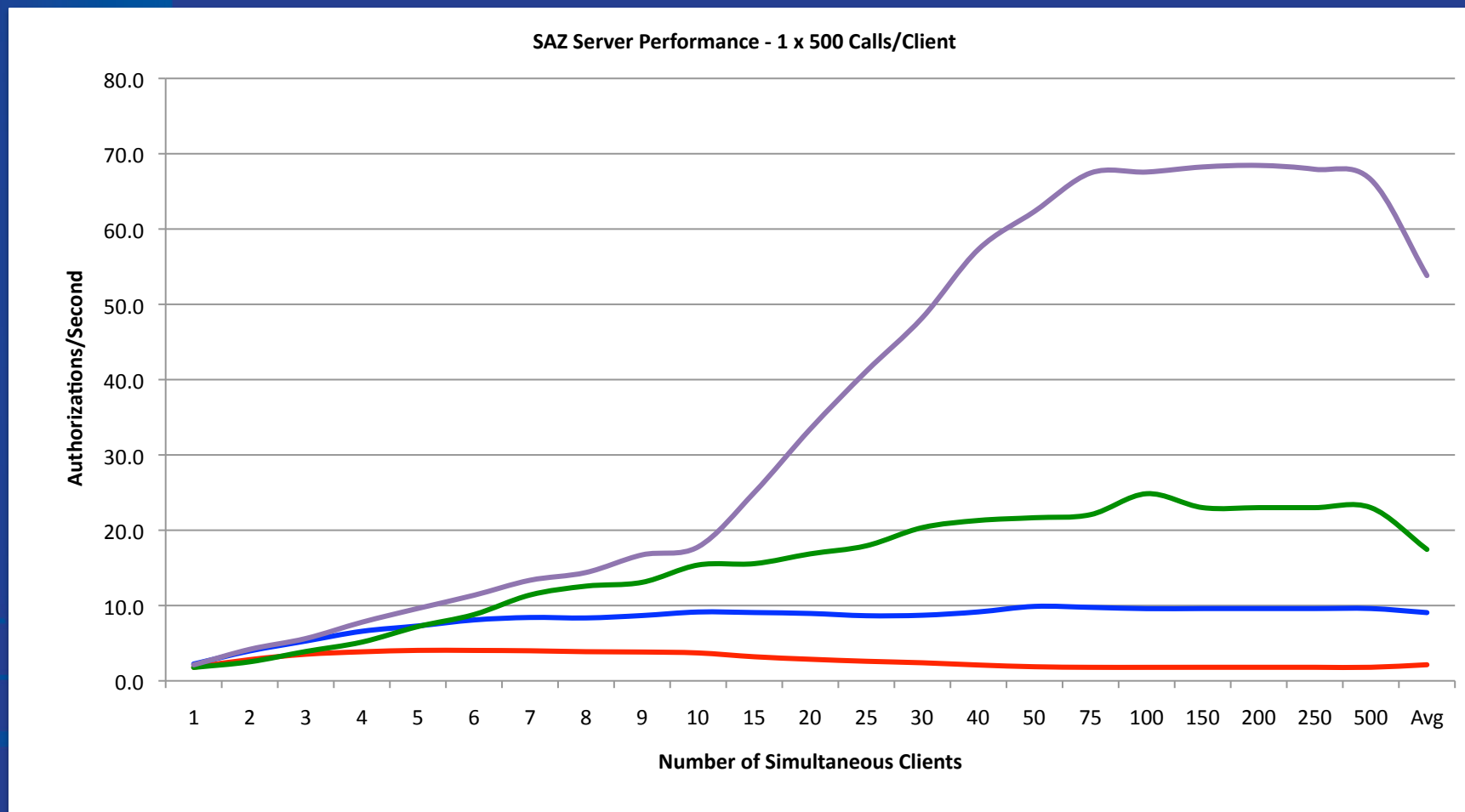
- Using development server
- Black = number of condor jobs,
- Red = number of SAZ network connections.
- Trigger @ 00:46:20,
- 25,000 Authorizations,
- 25,000 Success,
- 0 Failures,
- Complete @ 01:05:03,
- Elapsed time = 18m 43s,
- 22.26 Authorizations/sec.



# Comparison of Old & New Protocol



# Multiple Client SAZ production performance V2.0.1b (red) vs. V2.7.2 (purple)



# FermiGrid-HA2 Project

- Since the start of the FermiGrid project deployment in 2005, the physical machines were all in the FCC1 computer room – a single building with the corresponding power and network infrastructure:
  - Vulnerable to building issues (power issues 4X in past 18 months),
  - Vulnerable to network issues (6 in February 2011, 3 more in May 2011).
- The FermiGrid-HA high availability infrastructure had been designed to allow us to extend the design to support redundant services distributed across multiple buildings.
- The goal of the FermiGrid-HA2 project (started in March 2010) was to spread the systems and services between two buildings to lessen the chance of network cut or building outage disrupting all service.
  - Tuesday 24-May-2011 "Build & Test" – Move FermiGrid-HA2 Rack #1 to FCC2.
  - Tuesday 07-Jun-2011 "Go Live" – Move FermiGrid-HA2 Rack #2 to GCC-B.
  - <http://cd-docdb.fnal.gov/cgi-bin/ShowDocument?docid=3739>

# FermiGrid-HA2 Rack Layouts

- Two (almost) identically configured racks:
  - One with 120VAC power distribution,
  - One with 208VAC power distribution,
- Both currently located in FCC1 computer room,
- Currently waiting on network configuration changes,
- 1<sup>st</sup> rack will be moved to FCC2 computer room,
- 2<sup>nd</sup> rack will be moved to GCCB computer room.

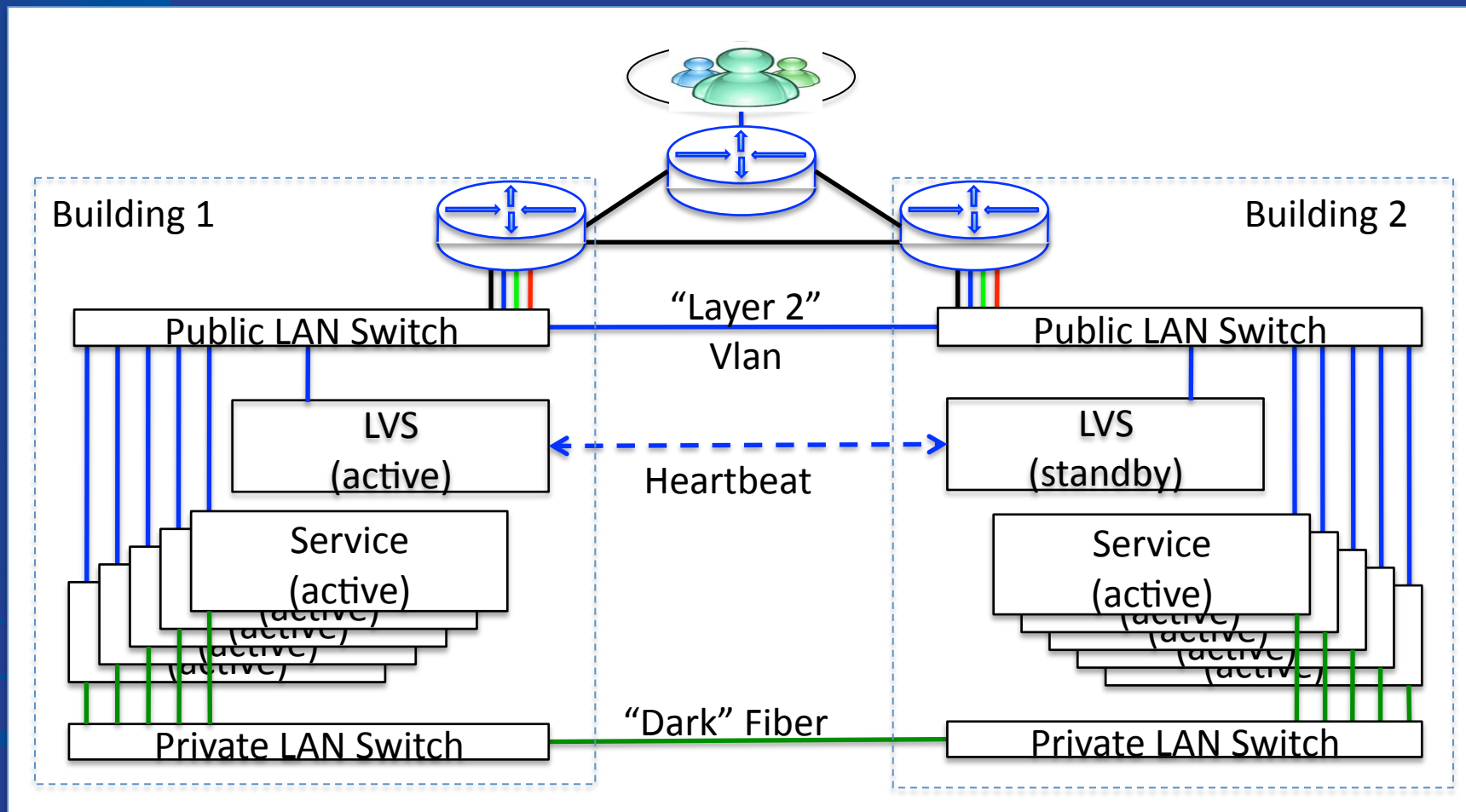
FermiGrid-HA2 Rack Front	Rack "U"	FermiGrid-HA2 Rack Rear
blank - 2U	42	Cisco Nexus / 2960G Public LAN Switch
fnpcsrv8 / blank - 1U	41	Cyclades AlterPath Console Server 16
fnpcsrv5 / fnpcsrv9	40	fnpcsrv8 / blank - 1U
fnpcsrv3 / fnpcsrv4	39	fnpcsrv5 / fnpcsrv9
blank - 1U	38	fnpcsrv3 / fnpcsrv4
d0osgsrv1 / d0osgsrv2	37	blank - 1U
blank - 2U	36	d0osgsrv1 / d0osgsrv2
fcdfsrv3 / fcdfsrv4	35	blank - 2U
fcdfsrv1 / fcdfsrv2	34	fcdfsrv3 / fcdfsrv4
fcdfsrv0 / fcdfsrv5	33	fcdfsrv1 / fcdfsrv2
blank - 2U	32	fcdfsrv0 / fcdfsrv5
Display / Keyboard / Mouse	31	Cyclades PM10-L30A 120VAC
Raritan MasterConsole MCCAT116 KVM	30	Cyclades PM10-L30A 120VAC
blank - 1U	29	Display / Keyboard / Mouse
Slave KDC - Sun Netra X1	28	Omniview PS3 16 port KVM
blank - 1U	27	Linksys SR2024 Private LAN Switch
gratia12 (FCC) / gratia13 (GCC)	26	APC Transfer Switch
gratia10 (FCC) / gratia11 (GCC)	25	blank - 1U
blank - 1U	24	gratia12 (FCC) / gratia13 (GCC)
ress01 / ress02	23	gratia10 (FCC) / gratia11 (GCC)
fermigrd5 / fermigrd6	22	blank - 1U
fermigrd2 / fermigrd3	21	ress01 / ress02
fermigrd1 / fermigrd4	20	fermigrd5 / fermigrd6
fermigrd0 / fermigrd7	19	fermigrd2 / fermigrd3
blank - 2U	18	fermigrd1 / fermigrd4
	17	fermigrd0 / fermigrd7
	16	Cyclades PM10-L30A 120VAC
	15	Cyclades PM10-L30A 120VAC
	14	
	13	
	12	
	11	
	10	
	9	
	8	
	7	
	6	
	5	
	4	
	3	
	2	
	1	



# FermiGrid-HA2 Pictures

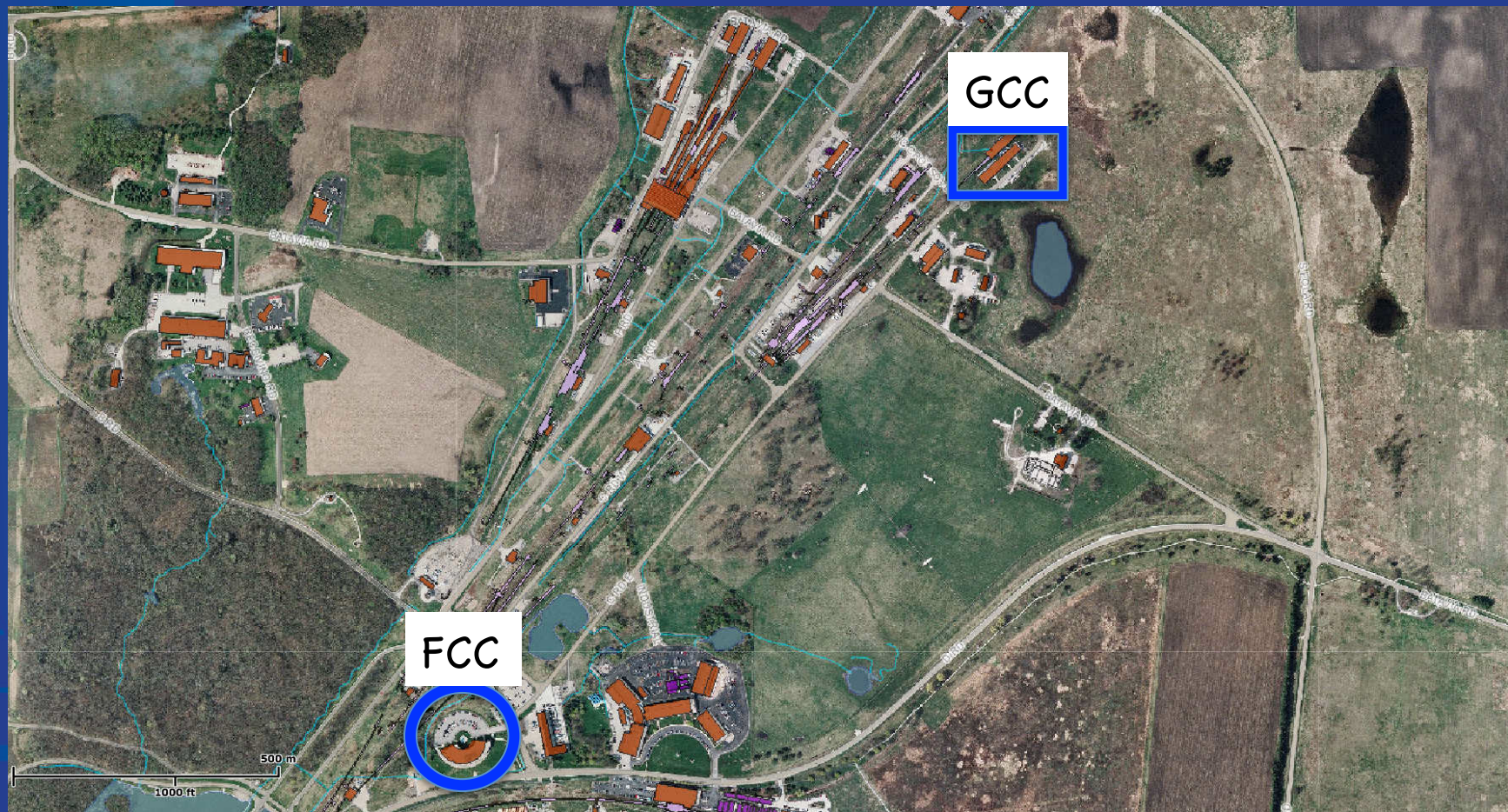


# FermiGrid-HA2 Network





# Geographical Redundancy



# FermiGrid-HA2 Project Results

- The FermiGrid-HA2 physical reorganization was completed on Tuesday 07-Jun-2011 at ~1300.
  - Critical services are now hosted in two data centers (FCC2 & GCC-B).
  - Non-critical services are split across the two data centers.
- The plan had been to utilize an upcoming scheduled power outage on 13-Aug-2011 for FCC2 to serve as the final acceptance test for the FermiGrid-HA2 project.
- At 1500 on Tuesday 07-Jun-2011, there was a cooling outage in the GCC-B data center. As a consequence, all systems in GCC-B were immediately shutdown by the facilities personnel tripping the main electrical panel breakers.
- FermiGrid-HA2 functioned exactly as designed.
  - The critical services failed back to the single remaining copy of the service on FCC2.
  - The non-critical services went to reduced capacity.
  - When power was restored at 1700, the second copy of the critical services transparently rejoined the service “pool”, and the non-critical services resumed operation.



# FermiGrid-HA and FermiGrid-HA2 (in automobile terms)

- Mazda Miata:
  - Affordable,
  - Performance,
  - Fun to drive,
  - Easy to maintain!



# GridWorks

## OSG Storage Test Stand

- Hardware acquired in 2009:
  - 8 x 3.0 GHz Intel Xeon, 16 GB memory, 1.5 TB disk.
- 5 physical systems:
  - gw014-gw019,
  - SLF5.4,
  - Virtualized using KVM,
  - 4 virtual machines per physical system,
  - Statically deployed cloud computing.

# FermiCloud

- A project to evaluate the technology, make the requirements, foster the necessary software development, and deploy the facility.
- Infrastructure-as-a-service facility:
  - Clients (developers, integrators, testers, etc.) get access to virtual machines without system administrator intervention.
  - Virtual machines are created by users and destroyed by the clients when no longer needed. (Idle VM detection coming in phase 2).
  - Testbed to let us try out new storage applications for the Grid and cloud.
- A private cloud:
  - On-site access only for registered Fermilab clients.
  - Can be evolved into a community or hybrid cloud with connections to Magellan, Amazon or other cloud providers in the future (AKA cloud bursting).
- Unique use case for cloud:
  - On the public production network,
  - Integrated with the rest of the Fermilab infrastructure.

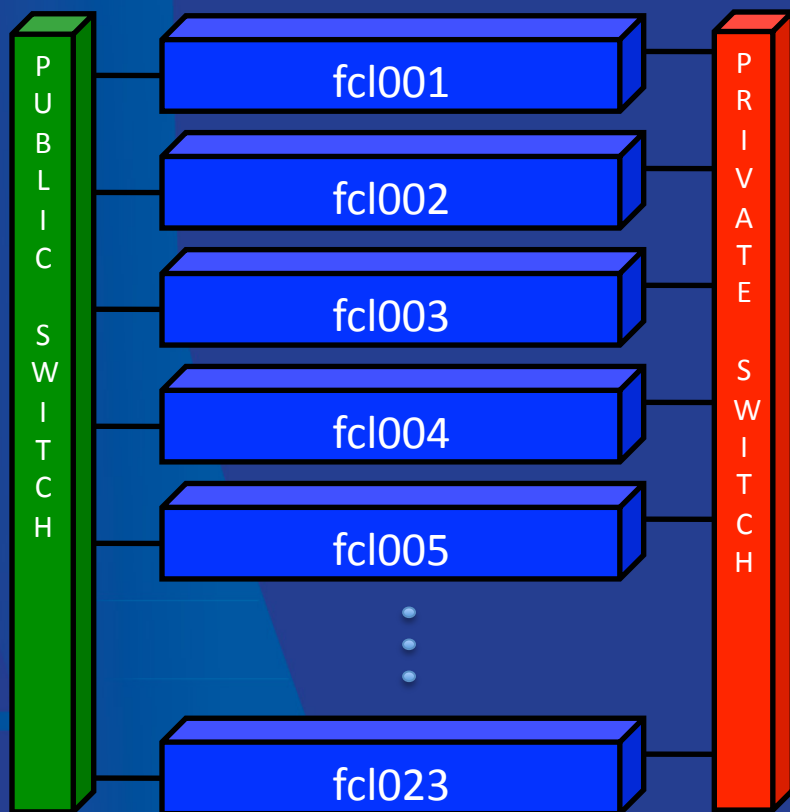
# FermiCloud Hardware



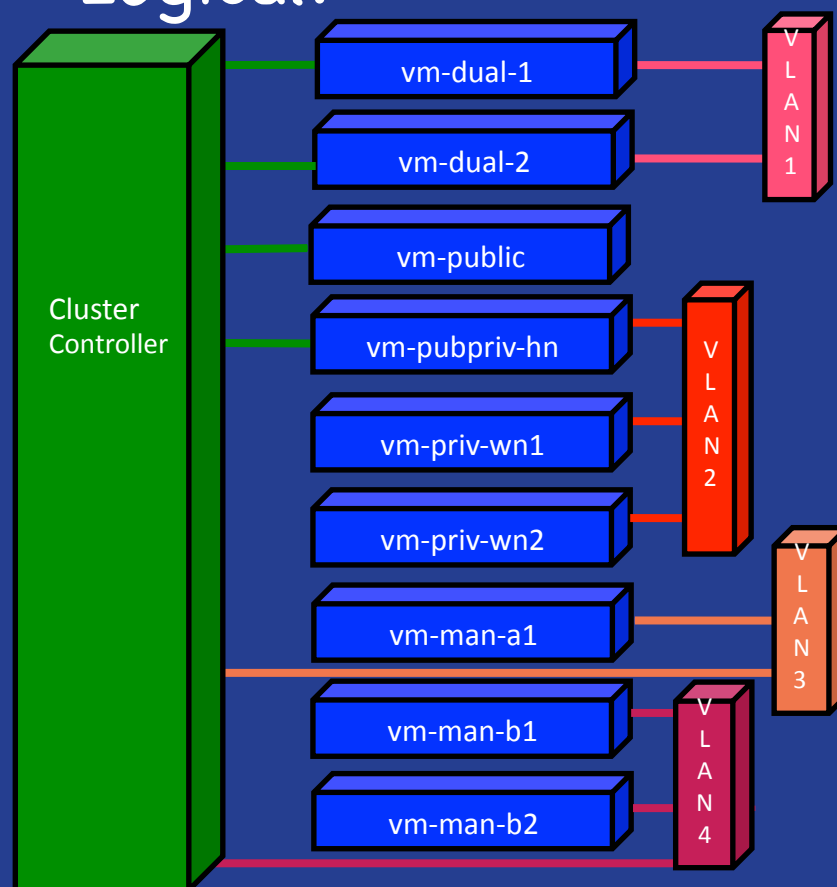
- Acquired in May 2010.
- Currently 23 systems located in GCC-B.
- Individual System:
  - 2 x Quad Core Intel Xeon E5640 CPU
  - 24GB RAM
  - Storage:
    - 2 x 300 GB SAS 15K rpm system disk.
    - 6 x 2TB SATA disk.
    - LSI 1078 RAID controller.
  - Connect-X DDR Infiniband
- Will likely buy SAN to attach to existing systems and split systems across two buildings later this year.

# FermiCloud Network Topology

- Physical:



- Logical:



# FermiCloud Project – Phase 1 (largely complete)

- Acquisition of FermiCloud hardware (done).
- Development of requirements based on stakeholder inputs (done).
- Review of how well open source cloud computing frameworks (Eucalyptus, OpenNebula, Nimbus) match the requirements (done).
  - The winner was OpenNebula (Nimbus was a close second).
- Storage evaluation (in process).
  - Lustre evaluation has been placed in CD-DocDB.
  - Hadoop, BlueArc evaluation nearly final,
  - OrangeFS yet to do.
- FermiCloud is being used by several Stakeholders today:
  - Grid & Cloud Computing Department, CD, CET, DMS, REX, CDF, DO, IF, CMS & OSG,



# FermiCloud Project – Phase 2 (underway today)

- Develop and deploy the necessary components to meet the requirements for selected cloud computing frameworks (Focus on OpenNebula and to a lesser extent Nimbus):
  - Gratia accounting, logging, monitoring, authorization, cloud bursting, etc.
  - Ted Hesselroth has delivered a pluggable authorization module for OpenNebula that works with DOEgrids and Fermilab KCA certificates.
- Vibrant collaboration with the open source Cloud Computing communities.
  - The pluggable authorization module for OpenNebula has been submitted back to the OpenNebula project and is in the process of being incorporated into the “trunk”.
- Improve the utilization of FermiCloud resources:
  - Reap (shut down and “shelve”) idle cloud VMs,
  - Boot up a “worker node” VM image to join an existing Grid cluster (GP Grid),
  - Allow utilization of otherwise idle resources.
- Develop and deploy cloud appropriate authorization, monitoring, accounting, logging, etc.
  - Will develop a cloud computing probe for Gratia accounting system.
- Production operation of various “small impact” services,
- Plan FermiCloud-HA and start acquisition of the necessary hardware.

# FermiCloud Project – Phase 3 (actively being developed)

- Production operation of services,
- Onboard additional production services & scientific stakeholders,
- Formal ITIL SLAs (Service Level Agreements):
  - “Development/Integration”,
  - Guaranteed Availability.
- Incremental deployment of FermiCloud-HA,
- Extend FermiCloud to support hybrid “cloud bursting”:
  - Run jobs on other Cloud providers (such as Amazon EC2).
- Running user/OSG provided virtual machines?

# FermiCloud Computing (In automobile terms)

- Resources on demand.
- When you need them for as long as you need them.



- Only “pay” for the resources you use.
- Minimize your resource usage.

# Grid & Cloud Computing Inventory as of April 2011

Mission	# Systems	# VMs	Virtualization Type
FermiGrid Production Services	8	52	Xen
CDF Grid Cluster Head Nodes	6	17	Xen
D0 Grid Cluster Head Nodes	2	7	Xen
GP Grid Cluster Head Nodes	7	33	Xen
Gratia Production	4	22	Xen
ReSS	2	10	Xen
FermiGrid Development	4	29	Xen
Gratia Development	5	30	Xen
CDF Test / Sleeper Pool	3	9	Xen
FermiGrid ITB	10	45	Xen
GridWorks	5	20	KVM
FermiCloud	23	66	KVM (with a little bit of Xen)
<b>Grand Total</b>	<b>79</b>	<b>340</b>	

**Average number of VMs per physical system: 4.30**

# Other Fermilab “Enterprise” Virtualization Projects

- Virtualization in SCF/FEF – General Physics Compute Facility (GPCF):
  - PaaS Facility,
  - Deployment of experiment-specific virtual machines for Intensity Frontier experiments,
  - Oracle VM (Commercialized RHEL+Xen).
  - High Availability configurations possible (presently limited by systems only in one building).
- Virtual Services Group:
  - Virtualization of Fermilab development and production core computing systems using commercial VMWare,
  - Windows, RHEL, SLF.
  - High Availability configurations possible.

# FEF VM Clusters

- 2 x Virtual Iron clusters
  - CDF Online webserver
  - D0 Offline
- 2 x Oracle VM clusters
  - CDF
  - GPCF

# What is GPCF?

- General Physics Computing Facility
- GPCF was created to solve a problem. We wanted to provide new and small experiments with inexpensive computing resources quickly.
- Additionally, GPCF allows us to consolidate moderately loaded one-off servers.





GPCF is not A fine Italian Sports CAR

It's not the cloud and we don't have Hadoop or Lustre...





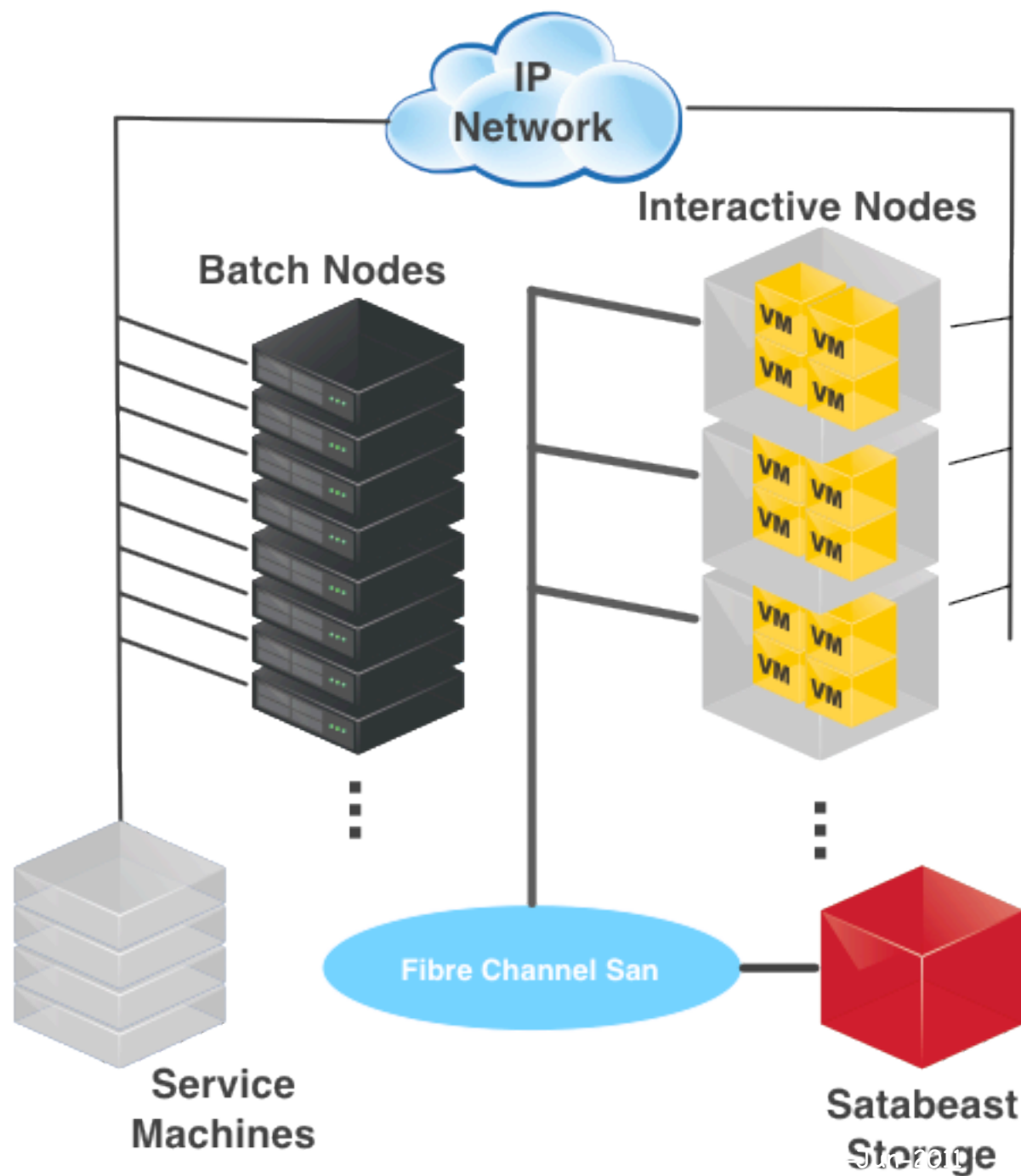
**GPCF is like A reliable MINIVan**

It will get you to your destination safely and in relative comfort.

# GPCF Info

- In production for almost a year!
- Used by 9 Intensity Frontier experiments to get real work done
- 13 interactive nodes running Oracle VM (~30 VMs)
- 10 x bare metal batch machines (Jobs can be submitted from any VM to run on batch nodes)
- 6 x bare metal service nodes for resource intensive work, e.g. staging data to be written to Enstore

# GPCF Topology



# GPCF Summary

- FY11 Upgrades
  - Additional 6 interactive nodes
  - Upgrade amount of memory in interactive nodes to 48GB
  - Reevaluate backend storage
- GPCF is a stable computing environment that gives emerging experiments an easy to use and flexible platform for getting work done.

# Virtual Services

# Virtual Services

- Created a general virtual infrastructure in 2010.
- Separate infrastructures based on VMware were run in TD, CD, and MIS since 2005.
- VMware vSphere is deployed on all host servers
- Target servers and desktops running Windows, RHEL, and SLF.
- Host servers and storage will be located in multiple data centers (FCC2, FCC3, ANL, and possibly GCC, WH)
- Provide support in diagnosing guest VM issues related to performance, configuration, and capacity.
- Provide monitoring and alerting for guest VM's.
- Provide assistance with P2V, V2V, and in the future V2P activities.
- By Q3 FY11, we will have 9 hosts running ~200 virtual machines.



# Virtual Services Server Hardware

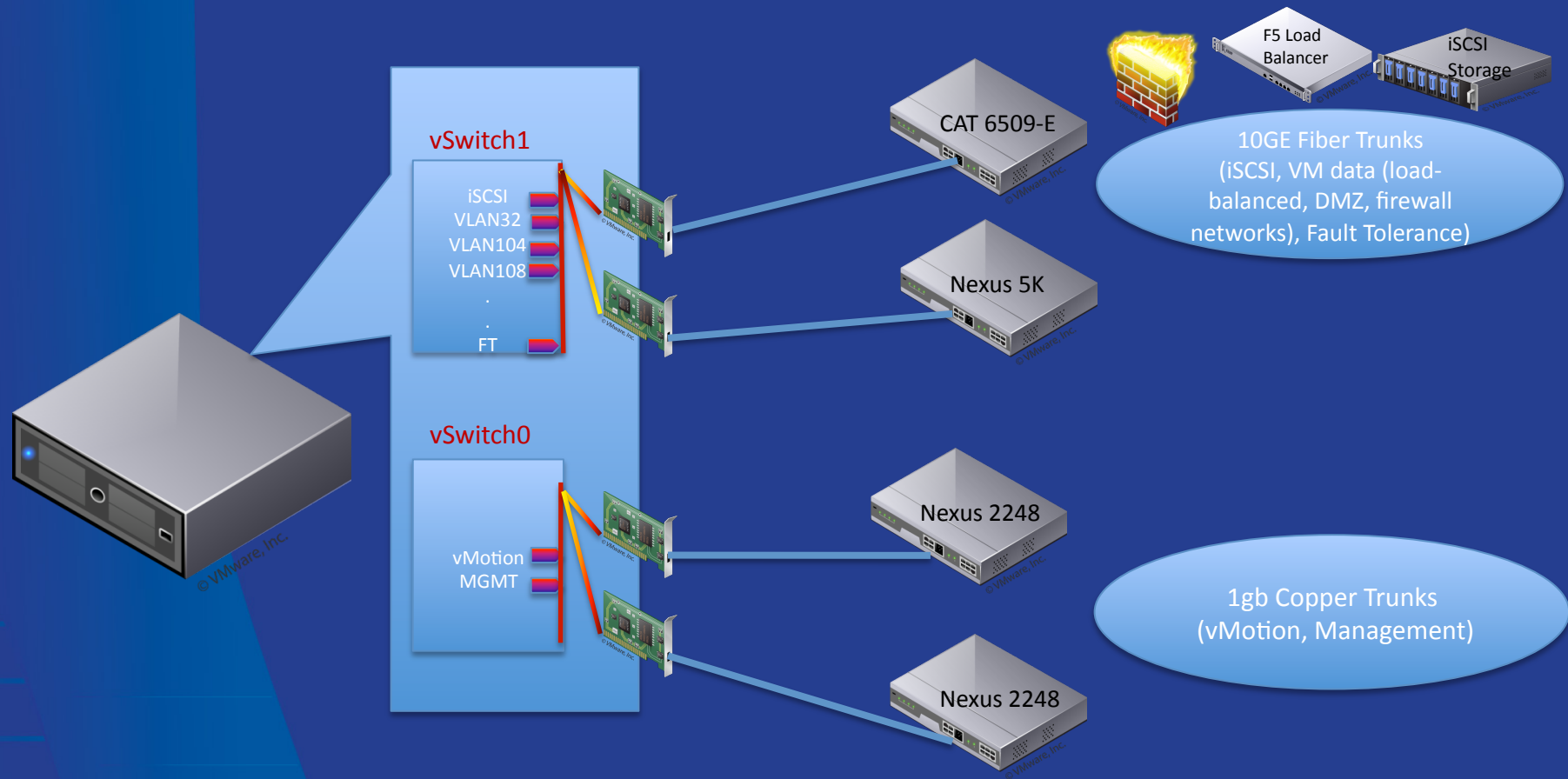
Deploying large capacity servers to minimize infrastructure overhead, such as SAN ports, network ports, rack space, power/cooling

- Currently 13 systems in 3 clusters (7 in FCC2, 2 in WH, 4 in TD's ICB).
- System Configurations:

- CD General Cluster individual system:
  - 128GB RAM
  - 4 x 6 core Intel Xeon 7460 CPU's (total=24 cores each)
  - Redundant dual port 10GE NIC's
  - Redundant dual port 8 Gb HBA's
- MIS Cluster individual system: (Retire 2 older servers in Q3)
  - 2@96GB RAM and 2@ 32GB RAM
  - 2@4 dual core AMD CPU's and 2@4 quad core AMD CPU's
  - 2-Quad 1GB NICs
  - Redundant dual port 8Gb HBA's
- TD Cluster individual system: (Consolidating into general cluster)
  - 48GB RAM
  - 2 x 4 core Intel CPU's
  - 2-Quad 1GB NICs (used for networking and iSCSI storage)
- Totals:
  - Host Servers: 13
  - CPU cores: 200
  - RAM: ~ 1TB
  - VM's and Templates: 160 (as of 2-4-2011)



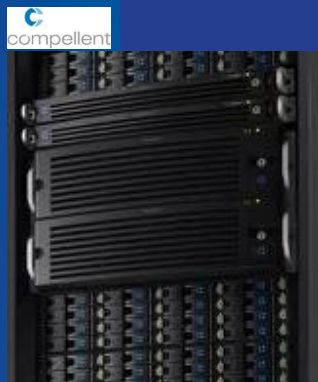
# Virtual Services Network Topology (CD general cluster)





# Virtual Services Fibre Channel Storage

- 2 – Compellent Storage Arrays
  - ~64TB total comprised of 32x1.0TB 7.2K RPM SATA disks, 68x450GB 15K FC disks, 6x146GB SSD's.
  - Key features include automated and policy-based tiered storage(SSD→FC→SATA), advanced thin provisioning, continuous data protection (snapshots, replays), thin replication, dynamic storage migration (allows you to migrate live data from one array to another on the fly).
  - Fault tolerance capabilities through dual paths from servers to disk drives, fully redundant power supplies and fans, and clustered controllers.
  - Used for DEV, QA, and PRD virtual machine instances.



# Virtual Services iSCSI Storage



- 4 x Dell EqualLogic PS6000E iSCSI SAN Arrays
  - 16TB each(16x1TB 7.2K SATA-II)
  - Comes with all backup and protection capabilities built in (including snapshots, clones, replication, and scheduling)
  - Fault tolerance capabilities through fully redundant and hot swappable components – standard dual-controllers, dual fan trays, dual power supplies and disk drives with hot spares.
  - Used for Test and DEV VM's as well as virtual desktops.
- 1 – Dell PowerVault MD3000i iSCSI SAN Array
  - 15TB (15x1TB 7.2K Sata)
  - Used for quick VM image backup and restores.
  - Entry level storage, but can take snapshots and replicate data with our current license.

# Virtual Services Stakeholders

- Finance Section/BSS (Domino, file servers)
- Directorate (Budget office, Audit, VMS)
- TD (Web, App, File server)
- Services/Projects (Web Servers, FTL, TeamCenter, Print Servers, Sharepoint, Teammate, Node Registration, Indico, Crystal Reports, DB Servers, MRTG, Meeting Maker, Plone, NIMI/Tissue, FIdM Dev, Wireless Control Server)

# Some Performance Measurements

# Performance Measurements/Limitations - 1

- Xen VM I/O performance measurements to BlueArc:
  - Xen read performance from BlueArc (~90 MB/sec) is comparable to “bare iron” (~100 MB/sec),
  - Xen write performance to BlueArc (~100 MB/sec) is comparable to “bare iron” (~100 MB/sec),
- KVM VM I/O performance measurements to BlueArc:
  - KVM read performance from BlueArc is comparable to “bare iron”,
  - KVM write performance to BlueArc is comparable to “bare iron”,
- Lustre VM I/O performance measurements:
  - Read performance from KVM Virtualized Lustre is comparable to “bare iron”,
  - Write performance to KVM Virtualized Lustre (~20 MB/sec) is significantly less than “bare iron” (~80 MB/sec).

# Performance Measurements/Limitations – 2

- Network equipment may affect performance measurements:
  - Cisco 6509 & 2960G – “Bare iron” read and write performance measured ~100 MB/sec.
  - Cisco 2248 – Initial “Bare iron” performance measurements were write ~100 MB/sec, read only 5-10 MB/sec.
    - Apparently caused by packet drops when the 10 Gb/s input “flow” had to be throttled into the 1 Gb/s switch port, the remaining 9 Gb/s was being dropped by the fabric extender.
  - Cisco 2248 – “Bare iron” read performance now much better ~80 MB/sec (changed configuration of 2248 to use “nodrop” queue).
- MySQL server performance:
  - A MySQL server deployed on a KVM VM responds to a simple query test with a system load in excess of 30.
  - An identically configured MySQL server deployed under Xen responds to the identical simple query test with a system load of under 1, which is ~equivalent to “bare iron”.

# Summary and Final Thoughts

# Workload

- Virtualization can generally deliver performance that is comparable to “bare iron”,
- There are cases where virtualization can actually deliver performance in excess of “bare iron”:
  - Non Uniform Memory Access (NUMA) Systems,
  - Use NUMACTL to lock virtual machine to a particular processor & associated memory,
  - Recent processors from Intel and AMD are NUMA.
- Not all workloads may be appropriate for Virtualization/Cloud Computing,
  - Example – a workload that requires all the resources of a system to accomplish its tasks,
  - Might still choose to virtualize this workload to aid in system maintenance or migration.



# Open Source vs. Commercial

- FermiGrid / GridWorks / FermiCloud use Open Source:
  - Xen and KVM hypervisors,
  - OpenNebula (and Nimbus),
  - Guest VMs are RHEL/SLF (and Windows in the future),
- GPCF uses commercial OracleVM (formerly Virtual Iron):
  - License is free, support is \$600 per year per system,
  - Xen hypervisor,
  - RHEL/SLF,
- Virtual Services uses commercial VMware:
  - License \$3K per CPU socket (processor),
  - Windows,
  - Can also run RHEL/SLF.

# Xen vs. KVM

- At the moment, some benchmarks and workloads show that Xen virtualization has better performance than KVM virtualization.
- Xen and KVM allow overbooking of CPU.
- KVM allows overbooking of RAM:
  - Can provision many more machines on cloud resources.
  - Kernel will share read-only pages and perform copy-on-write.
- KVM is incorporated into the “stock” SL(F,C) Kernel as of SL(F) 5.4+, so it is significantly easier to virtualize a “bare iron” machine.
- However...
  - While the SL(F,C) upstream vendor still lists Xen as the default hypervisor in the 5.x distribution, they have announced that KVM is the future.
  - FermiGrid has observed that support for Xen based virtualization is declining in the SL(F,C) upstream vendor distribution (time synchronization issues with a 64 bit hypervisor and 32 bit virtual machine hardware clocks).

# Security

- Virtualization and Cloud Computing do not eliminate security issues.
- Virtualization actually offers an additional “surface” to potentially attack – the virtualization layer.
- Cloud computing, while effective in delivering resources on demand, can also increase your risk (as shown by the recent Amazon EC2/EBS outage and data loss).

# Final Words

- Virtualization can deliver:
  - Flexibility, Availability, & Performance.
- Virtualization is not the full solution:
  - A very good tool to have in your toolbox,
  - Must consider the entire life-cycle,
  - Additional mitigations may be necessary.
- Workloads that are appropriate for Virtualization/Cloud Computing must be carefully deployed to insure adequate performance.
  - Memory bandwidth & utilization,
  - Local and remote file systems,
  - Network utilization,
  - CPU.
- Testing and Monitoring are essential:
  - Testing at scale is the key to reliable operations at scale
  - Monitoring is the key to assuring service delivery

# Virtualization alone will not address this vulnerability:

- Wildfire destroys half of town of 9,800:



Fin

Questions?